

# Interroger un corpus par le sens

## Une approche linguistique

Bernard JACQUEMIN  
Institut des Sciences Cognitives  
CNRS – UMR5015  
67, Bd Pinel  
69675 Bron cedex (France)  
Bernard.Jacquemin@isc.cnrs.fr

### Résumé

Constatant que les méthodes statistiques dominantes en traitement de l'information ne peuvent résoudre certaines difficultés, je propose une approche centrée sur les méthodes linguistiques symboliques afin d'identifier la contribution que ces dernières peuvent apporter au domaine. Elle s'appuie sur une identification du sens des mots et des relations entre ces mots pour proposer des reformulations d'énoncés sans changement de signification. Les reformulations sont parasynonymiques et dérivationnelles et permettent de trouver une information textuelle quelle que soit la formulation de l'information ou de la requête.

**Mots-clefs** : Sémantique lexicale; désambiguïsation; question-réponse; extraction d'information; corpus; génération; reformulation; analyse syntaxico-sémantique; synonymie; dérivation morphologique; dictionnaire électronique

### Abstract

In textual knowledge management, statistical methods prevail. Nonetheless, some difficulties cannot be overcome by these methodologies. I propose a symbolic approach using a complete textual analysis to identify which analysis level can improve the the answers provided by a system. The approach identifies word senses and relation between words and generates as many rephrasings as possible. Using synonyms and derivative, the system provides new utterances without changing the original meaning of the sentences. Such a way, an information can be retrieved whatever the question or answer's wording may be.

**Keywords**: Lexical semantics; Word Sense Disambiguation; question answering; information extraction; corpus; generation; rephrasing; parsing; semantic analysis; synonymy; derivational morphology; electronic dictionary

## 1 Introduction

La société actuelle a fait de la maîtrise de l'information un enjeu de savoir autant que de pouvoir. Cependant, face à la profusion des sources d'information, face à l'enchevêtrement ingérable des données elles-mêmes, personne n'est plus capable de fournir un accès rapide à un élément d'information précis. Les initiatives qui visent à élaborer une méthode automatique de gestion de l'information capable d'ordonner des masses de données sont dès lors bienvenues.

Dans les champs de recherche liés à l'information textuelle électronique, et notamment la tâche de question-réponse, plusieurs méthodes ont vu le jour, qui permettent de confronter les données de la question avec celles contenues dans un texte. Si les données correspondent, on considère que la réponse à la question posée est dans le contexte immédiat de l'information commune à la question et à la bribe de texte. Diverses conférences internationales ont également vu le jour, dont l'objet est l'évaluation des systèmes proposés : TREC, CLEF, NTCIR. . .

Il reste toutefois que toutes ces méthodes fonctionnent sur base d'un même schéma. Il s'agit en effet d'appréhender la question, de l'analyser pour la débarrasser de tout élément perturbateur, et d'en effectuer une expansion destinée à contrer les variations de forme qui peuvent se présenter dans les documents interrogés. De plus, les meilleurs systèmes plaident tous pour des approches capable de gérer au mieux le caractère langagier des textes. Par exemple, [Hull, 1999] proposait déjà l'exploitation des résultats d'une analyse morphologique pour indexer les éléments significatifs tant dans les requêtes que dans les réponses. Par la suite, [Ferret *et al.*, 2002] ont proposé une certaine intégration de la syntaxe, capable de reconnaître un certain nombre d'entités nommées ainsi que de déterminer la nature de la réponse à la question. Le meilleur système actuellement disponible [Harabagiu *et al.*, 2000] s'appuie sur des notions sémantiques issues du réseau WordNet, ainsi que sur un moteur d'inférences logiques pour proposer un niveau très élevé de réponses correctes.

De plus, une consultation, même rapide, des publications liées aux campagnes d'évaluation du domaine de question-réponse [Voorhees et Buckland, 2004, Peters *et al.*, 2005, Kando et Ishikawa, 2005] permet de constater que les qualités qui distinguent le résultat des différentes approches résident dans leur capacité à mieux appréhender la langue, à obtenir une meilleure analyse linguistique non seulement de la question, mais également des réponses possibles. Malgré les vertus reconnues des modules d'analyse linguistique, il est pourtant étrange de constater qu'ils n'occupent qu'une place générique dans toutes ces approches à dominante statistique, et qu'aucune recherche n'est actuellement menée pour leur accorder un statut plus central qui pourrait encore améliorer le fonctionnement des logiciels.

Cet article s'appuie sur les constats précédents pour proposer une méthode où l'analyse linguistique occupe une place prépondérante à tous les niveaux du système. Nous allons d'abord présenter les contraintes propres à une analyse morpho-syntaxique et sémantique d'énoncés textuels et exposer les choix auxquels elles nous ont amené. Ensuite, nous présenterons brièvement les outils d'analyse que nous avons utilisés dans notre approche, ainsi que les ressources lexico-sémantiques que nous avons exploitées, et les adaptations qui ont été nécessaires. Après cela, nous présenterons la construction d'une structure informationnelle qui fournit un accès à chaque élément d'information contenu dans une base textuelle. Enfin, nous présenterons la méthode permettant de trouver la réponse à une question posée en français dans cette base textuelle. Finalement, nous présenterons quelques perspectives futures pour ce sujet de recherche, et notamment l'exploitation de cet outil en manipulation linguistique de corpus textuels.

## 2 Une méthode linguistique de structuration textuelle

L'examen des systèmes de question-réponse existants nous a donc amené à élaborer une stratégie bien différente, centrée sur une analyse linguistique des énoncés. Cette démarche s'inscrit toutefois dans la tradition du domaine, puisque la définition de l'information considérée appartient à une **perspective** syntaxico-sémantique et lexico-sémantique. Il s'agit en

effet d'identifier les unités lexicales porteuses de sens, considérées comme l'information élémentaire, les relations (syntaxiques) entre ces éléments, ainsi que le sens lexical des lexèmes en contexte. D'autre part, les variations dans la formulation d'une même information sont classiquement compensées par une expansion de l'énoncé, où synonymes, hypéronymes, holonymes et autres dérivés morphologiques interviennent à plaisir.

L'expansion d'énoncé pratiquée généralement dans le domaine de question-réponse s'applique à la requête. Le principe en est simple : à chaque unité lexicale considérée comme significative est associée une liste de mots qui lui sont jugés équivalents, sous forme de synonymes, d'hyponymes et hypéronymes, de dérivés, etc. Ces listes d'expansions sont utilisées disjonctivement au lexème original lorsque la requête est proposée à un moteur de recherche. Mais si une telle expansion permet effectivement de résoudre dans un grand nombre de cas les problèmes de formulation, seule la maîtrise du sens de l'énoncé à expander permet de sélectionner les reformulations qui conviennent dans le contexte courant. La figure 1 page 3 distingue les expansions correctes (A) des expansions erronées (B), susceptibles d'apporter des réponses inadéquates.

Question : « De quel **chef** Domitien est-il le **successeur** ? »

A	général	héritier
	empereur	succéder
	...	...
B	cuisinier	remplaçant
	cheveu	succédané
	...	...

Réponse :

Second fils de Vespasien, Domitien succéda à l'empereur Titus et poursuivit la remise en ordre de l'État.

FIG. 1 – Exemple d'expansion d'énoncé : le problème du sens

Une sélection du sens des unités lexicales de la requête est donc nécessaire pour que l'expansion puisse être effectuée en fonction du sens original, afin de limiter le bruit. Cependant, on sait depuis [Weaver, 1949] toute l'importance que prend le contexte – et même son contexte syntaxique [Reifler, 1955] – dans le choix du sens lexical d'un mot dans une phrase (désambiguïsation sémantique lexicale). De plus, le simple bon sens permet de constater que les questions que l'on peut poser au système sont généralement plus courtes que des phrases rédigées dans un document. Le contexte y est donc moins important que dans les textes interrogés. Par ailleurs, les grammaires syntaxiques des outils d'analyse existants fonctionnent habituellement moins bien sur des phrases interrogatives que sur des énoncés affirmatifs. Il est donc moins efficace de traiter la sémantique lexicale d'un lexème s'il apparaît dans une requête que s'il figure dans la base textuelle elle-même.

Dès lors, une démarche qui applique prioritairement une analyse syntaxico-sémantique aux documents plutôt qu'à la requête s'est imposée. Ce choix est d'autant plus opportun qu'il s'intègre pleinement à l'indexation du contenu des documents, indispensable lors de la phase de recherche, et qui consiste à recenser le contenu des documents. De plus, une telle

approche centrée sur les documents présente comme avantage pratique de distinguer nettement et chronologiquement l'analyse et l'expansion des énoncés, et la phase d'interrogation. De la sorte, les traitements les plus lourds sont appliqués préalablement, et l'interrogation de la structure de l'information se fait de manière presque instantanée.

Le système comporte donc deux niveaux de fonctionnement. Le premier consiste à analyser les documents d'un point de vue morphologique, syntaxique et sémantique, puis à leur appliquer une expansion à l'aide d'informations provenant de ressources lexico-sémantiques (adjonction d'*enrichissements*), et à stocker les résultats dans des index constituant une structure de toute l'information textuelle [Roux et Jacquemin, 2002]. La seconde étape a pour objet l'interrogation de cette structure à l'aide de questions ordinaires.

### 3 Outils d'analyse textuelle

Comme nous l'avons indiqué plus haut, divers outils d'analyse interviennent dans l'identification des éléments d'information présents dans les documents. Il s'agit d'identifier les éléments eux-mêmes, c'est-à-dire les mots significatifs, au travers d'une analyse morphologique ; ensuite, les relations entre ces mots grâce à l'analyse syntaxique ; enfin, la désambiguïsation sémantique permet de connaître la signification des mots dans leur contexte d'apparition. Voici une rapide description de ces outils.

#### 3.1 L'analyseur morphologique

L'analyseur morphologique NTM (*Normalizer, Tokenizer, Morphological analyzer*) que nous avons utilisé est un système de transducteurs à états finis développé au Centre de Recherche européen de Xerox (XRCE) [Aït-Mokhtar, 1998]. Ce système prend en entrée n'importe quelle chaîne de caractères en français et y applique des traitements de normalisation, de segmentation s'il s'agit de plusieurs unités lexicales, et propose les différentes analyses morphologiques possibles pour chacun des segments identifiés. La figure 2 page 4 permet d'identifier les traitements appliqués par NTM à une phrase proposée en entrée. Il permet d'obtenir une version normalisée de chaque unité lexicale, son lemme et les informations morphologiques qui y sont associées sous la forme de traits attachés à la forme de départ.

Son deuxième fils [...]

Mot du texte	lemme	Analyse morphologique	Traits ajoutés
son	son	+PP3S+InvGen+SG+Poss	
son	son	+Masc+SG+Noun+	+SOM+AGR

PP3S	Pronom personnel 3ème sg	Masc	Masculin
InvGen	Invariable en genre	SG	Singulier
SG	Singulier	Noun	Nom
Poss	Possessif	SOM	Relatif au corps
		AGR	Agriculture

FIG. 2 – Exemple d'analyse morphologique par NTM

Cet analyseur présente également une qualité qui a pu être exploitée avec succès. En effet, sa conception sous forme de transducteurs, et la présentation de ses résultats sous forme de traits attachés aux unités lexicales, permettent d'ajouter aisément certaines informations lexicales qui peuvent être utiles pour les traitements ultérieurs. Ainsi, on peut voir dans l'exemple que nous avons ajouté aux lexiques existants des informations sémantiques extraites d'un dictionnaire, qui seront utilisées ultérieurement lors de la phase de désambiguïsation. Dans l'intervalle, cette information subsiste attachée aux unités lexicales, mais elle reste virtuelle dans la mesure où elle n'intervient ni dans la désambiguïsation catégorielle, ni dans l'analyse syntaxique.

### 3.2 L'analyseur syntaxique

L'analyseur syntaxique XIP (*Xerox Incremental Parser*) [Roux, 1999] est un moteur d'analyse syntaxique basé sur des grammaires de réécritures incrémentales. Il permet d'effectuer le cas échéant une désambiguïsation catégorielle d'énoncés étiquetés morphologiquement mais non désambiguïsés. Il propose surtout une analyse syntaxique de surface robuste de ces énoncés sous forme de dépendances entre des nœuds représentés sous la forme des unités lexicales équivalant à la tête des syntagmes minimaux (*chunks*) concernés. Une représentation en arbre de chaque phrase, ainsi qu'un découpage en syntagmes minimaux sont également proposés, mais ils ne sont pas utilisés ici.

Énoncé : « Il reconstruisit Rome ruinée par les incendies. »

Extraction des dépendances :

SUBJ(reconstruisit,Il)	2e argument sujet du 1er argument
SUBJ(ruinée,incendies)	
VMOD[INDIR](ruinée,par,incendies)	3e argument compl. agent 1er argument
VARG[DIR](reconstruisit,Rome)	2e argument COD du 1er argument
NMOD[ADJ](Rome,ruinée)	2e argument épithète du 1er argument

FIG. 3 – Exemple d'analyse syntaxique par XIP

La figure 3 page 5 permet d'évaluer les possibilités de XIP et d'illustrer son mode de représentation des relations syntaxiques par dépendances. On peut également voir le travail de certains traits, exclusivement syntaxique ici, et portant sur la nature des dépendances (DIR et INDIR respectivement sur les dépendances VMOD et VARG, ainsi que ADJ sur NMOD). XIP applique des règles contextuelles qui permettent d'évaluer des nœuds et des traits portant sur des nœuds appartenant à un même contexte pour construire des syntagmes minimaux et des dépendances. Ce mode de fonctionnement est très intéressant car s'il permet de travailler sur des indications lexico-morphologiques pour la désambiguïsation catégorielle et des données lexico-syntaxiques pour la construction des dépendances syntaxiques, il n'y a pas d'obstacle à son utilisation dans une perspective lexico-sémantique.

### 3.3 Le désambiguïsateur sémantique

Le système de désambiguïsation sémantique présenté de le cadre de cette étude est une évolution de la méthode de [Brun *et al.*, 2001], qui reposait sur l'exploitation de l'analyse

syntaxico-sémantique d'un dictionnaire utilisé comme corpus sémantiquement étiqueté. Le présent système [Jacquemin, 2003] exploite l'information du *Dictionnaire des verbes français* [Dubois et Dubois-Charlier, 1997] et de son complément des autres catégories grammaticales (ces deux dictionnaires complémentaires seront désormais désignés sous le nom générique *Dubois*). Ces dictionnaires répartissent l'information fournie non par unité lexicale, mais par sens de chaque unité lexicale. De la sorte, chaque information fournie par le dictionnaire est discriminante pour le sens concerné d'un mot donné.

Le fonctionnement du désambiguïsateur se fait en deux temps : d'abord l'analyse du dictionnaire avec création de règles conditionnelles de désambiguïsation sémantique, basées sur un schéma syntaxique, et ensuite l'application de ces règles à des mots en contexte, sur base des contextes syntaxico-sémantiques fournis d'abord par les étiquettes sémantiques ajoutées à NTM, ensuite par les dépendances issues de l'analyse syntaxique de XIP.

L'information qui peut être extraite du dictionnaire régit le type de règles qui peuvent être construites. Dans le cas du Dubois, l'information peut être diverse et se présenter sous forme purement syntaxique (p.ex. « Je bois » *vs* « Je bois de l'eau » avec l'indication de transitivité), syntaxico-sémantique (p.ex. « embrasser quelque chose » *vs* « embrasser quelqu'un » avec sous-catégorisation du complément direct), lexico-syntaxique avec l'analyse des exemples et la conservation des relations impliquant le mot considéré comme autant de schémas typiques (« le général remporte la victoire » implique la dépendance `VARG[DIR](remporter,victoire)`, avec le mot victoire comme complément direct de remporter) ou sémantico-syntaxique (généralisation de la dépendance extraite d'un exemple grâce aux traits sémantiques correspondant à une unité lexicale : `VARG[DIR](remporter,[MIL])`, où le trait MIL pour militaire est le trait sémantique de victoire, qu'il remplace).

Comme les règles de désambiguïsation doivent répondre au contexte syntaxique, comme le stipule [Reifler, 1955], et qu'elles sont conditionnelles – puisque la conformité d'un contexte à une information issue d'un dictionnaire implique la sélection du sens correspondant – elles répondent à toutes les conditions pour en faire une grammaire dans XIP. C'est donc à cette syntaxe que les règles de désambiguïsation doivent correspondre, ce qui évitera la création d'un moteur d'application des règles particulier. Les résultats d'une désambiguïsation sémantique apparaîtront donc comme des dépendances extraites par XIP ou comme des traits sur des dépendances ou des nœuds de XIP. La figure 4 page 7 permet de comprendre le mode de construction de règles de désambiguïsation sémantique à partir de l'information contenue dans le dictionnaire. Les données syntaxiques ou lexicales sont formalisées sous la forme d'une dépendance XIP et les données sémantiques sous la forme de traits sur les nœuds.

L'application des règles se fait au travers de l'analyse des énoncés par XIP. La mise en correspondance de l'analyse syntaxique de XIP avec le schéma syntaxico-sémantique d'une règle de désambiguïsation implique l'application de la règle, ou la sélection du sens correspondant à cette règle. Le sens sélectionné est indiqué sous la forme d'un trait associé avec l'unité lexicale considérée. Les informations sémantiques ajoutées par NTM servent à l'application des règles impliquant des indications sémantiques.

## 4 Adaptation des ressources lexicales

On a déjà pu voir que les données lexicales étaient capitales pour l'approche proposée ici. Elles le sont non seulement dans la perspective de l'analyse sémantique, où l'information

Exemple extrait du Dubois pour « *remporter* » au sens 03 *gagner* :

« Le général remporte la victoire ».

Dépendances extraites de l'exemple :

SUBJ(**remporter**, **général**)

VARG[DIR](**remporter**, **victoire**)

- Construction d'une règle lexico-syntaxique de désambiguïsation :

**remporter** : VARG[DIR](*remporter*, **victoire**)

⇒ remporter 03 « gagner »

→ apparition de victoire comme complément direct de remporter implique le sens 03 « gagner »

- Construction de la règle sémantico-syntaxique correspondante :

victoire → trait sémantique : MIL (militaire)

**remporter** : VARG[DIR](*remporter*, **MIL**)

⇒ remporter 03 « gagner »

→ apparition d'un mot comportant le trait MIL (militaire) comme complément direct de remporter implique le sens 03 « gagner »

FIG. 4 – Exemple de construction des règles de désambiguïsation

syntaxico-sémantique distribuée par sens des entrées est prépondérante, mais aussi dans une optique d'expansion d'énoncés. En effet, cette expansion est effectuée par remplacement d'unités lexicales originales par d'autres, qui peuvent leur être substituées avec un minimum de modifications de sens. Ce sont donc deux types de modifications lexicales qui sont réalisées : la synonymie, et la dérivation morphologique. Dans une certaine mesure, le dictionnaire Dubois est à même de fournir les indications permettant de procéder à ces transformations.

En effet, un des champs informationnels de ce dictionnaire de référence fournit des synonymes, tandis qu'un autre procure des indications relatives à la dérivation. Toutefois, les synonymes sont invariablement au nombre de deux, ce qui est généralement insuffisant pour couvrir l'ensemble des transformations synonymiques possibles. D'autre part, les indications de dérivations se basent sur une racine et des suffixes, qu'un système automatique est difficilement à même d'interpréter correctement. Dès lors, d'autres ressources et outils doivent être exploités pour combler les lacunes du Dubois.

#### 4.1 Adjonction et répartition de synonymes

Pour ajouter une information synonymique au Dubois, nous avons utilisé trois ressources lexico-sémantiques : EuroWordNet [Vossen, 1998, Catherin, 1999], le dictionnaire des synonymes de [Bailly et Toro, 1947], et un dictionnaire multilingue utilisé comme outil chez Memodata. Tous fournissent des synonymes, mais leur répartition par sens, quand elle existe, ne correspond pas à celle du Dubois. Il a donc fallu les redistribuer. Nous avons

élaboré une méthode qui le fait automatiquement, décrite ici [Jacquemin, 2004b].

Cette procédure établit pour chaque entrée de chaque dictionnaire la liste des synonymes sans faire de distinction entre les sens différents que cette entrée peut avoir. Ensuite, à chaque synonyme proposé, elle associe toutes les étiquettes sémantiques qui lui sont attachées dans le dictionnaire Dubois. Puis, pour l'entrée considérée, chaque sens du Dubois est considéré successivement : lorsqu'une des étiquettes sémantiques du synonyme proposé est identique à celle du sens courant de l'entrée, il est considéré comme un synonyme valable pour ce sens et ajouté au champ de synonymie du Dubois. La même opération est effectuée pour chaque entrée de chaque dictionnaire de synonymes, puis les doublons sont éliminés. La figure 5 page 8 illustre la procédure suivie.

ravir (sens n°2, « voler »)

Synonymes proposés :	<b>enlever</b>	étiquette sémantique	<b>SOC</b> / LOC / TEX...
	charmer		PSY / OCC
	...		...

Étiquette sémantique de ravir (2) : **SOC**

⇒ synonyme ajouté : **enlever**

FIG. 5 – Répartition des synonymes par sens du mot original

## 4.2 Génération de dérivés

D'autre part, l'information contenue dans le Dubois ne permet pas d'effectuer automatiquement la génération des formes dérivées à partir d'une vedette du dictionnaire. Par contre, cette information peut se révéler suffisante pour identifier une proposition de dérivation et confirmer sa validité. L'outil de dérivation morphologique proposé par [Gaussier, 1999] peut dès lors être utile puisqu'il permet de générer, pour un mot proposé, un très grand nombre de candidats dérivés qui sont également des lexèmes attestés dans le lexique, à condition de lui laisser un maximum de latitude en diminuant au maximum les contraintes de génération. Les données de dérivation indiquées dans les champs correspondants du Dubois permettent ensuite, par identification du suffixe et de certaines caractéristiques de la racine, de ne conserver pour chaque sens que les dérivés prescrits par le dictionnaire.

## 5 Construction de la structure informationnelle

Comme on l'a vu plus haut, la structure de l'information est constituée d'index comprenant l'ensemble des données contenues dans les dictionnaires, et permettant d'avoir accès directement à la bribe de texte considérée comme intéressante dans la base documentaire. Cette structure est constituée d'abord du résultat de l'analyse des textes, c'est-à-dire des unités lexicales identifiées lors de l'analyse morphologique, ainsi que des traits morphologiques qui y sont associés, puis des relations syntaxiques entre ces lexèmes, ainsi que des traits syntaxiques associés soit aux dépendances, soit aux unités lexicales, et enfin des traits sémantiques dénotés lors de la désambiguïsation sémantique (numéro de sens,



indications sémantiques propres à ce sens dans le Dubois), qui portent uniquement sur les unités lexicales. La figure 6 page 9 montre comment la structure informationnelle est construite à partir d'un énoncé et de son analyse : les traits sont représentés entre crochets, les dépendances en majuscules et les unités lexicales en minuscules.

« ... Domitien succéda à l'empereur Titus. ... »

```
SUBJ(succéda[sn=1], Domitien[proper])
VARG[INDIR](succéda[sn=1,], à, empereur[humain,])
NN(empereur, Titus[proper])
```

FIG. 6 – Construction du « squelette » de la structure informationnelle

Dans un deuxième temps, la structure informationnelle est enrichie par l'expansion des énoncés. Les synonymes sont ajoutés disjonctivement aux dépendances dans lesquelles apparaissent les unités lexicales originales. Par contre, les formes dérivées ne peuvent être placées de même dans la structure informationnelle, car elles appartiennent le plus souvent à une catégorie grammaticale différente de celle du lexème original dont elles sont dérivées, et ne présentent pas une construction syntaxique similaire. Pour conserver une signification aussi proche que possible de l'énoncé original, il s'agit donc de reformuler l'énoncé pour qu'il intègre la forme dérivée. Pour cela, nous avons étudié le processus de dérivation : pour chaque catégorie grammaticale originale, pour chaque catégorie grammaticale dérivée, pour chaque type suffixal de dérivation, nous avons sélectionné au hasard trois exemples de dérivation dans le Dubois, et nous avons observé de quelle manière on peut remplacer l'original par le dérivé dans vingt contextes réels obtenus sur le Web. À partir de là se sont dégagés des patrons de dérivation qui permettent à partir d'un contexte sémantico-lexico-syntaxique, d'identifier le schéma syntaxique d'apparition de l'original et d'inférer un schéma syntaxique de dérivation. La figure 7 page 10 permet de comprendre de quelle manière les différents enrichissements sont ajoutés à la structure originale pour constituer des expansions des énoncés originaux, que ce soit par synonymie ou dérivation.

## 6 Interrogation de la structure

L'interrogation de la structure informationnelle peut être effectuée de nombreuses façons, dans la mesure où il suffit d'effectuer une recherche sur un ou plusieurs éléments stockés dans les index pour obtenir instantanément les énoncés d'apparition de ces éléments. Dans le cadre de l'application de question-réponse, c'est à une question en langue naturelle que le système doit apporter une réponse [Jacquemin, 2004a]. L'information contenue dans la question doit donc être convertie dans un format similaire à celui de la structure informationnelle, c'est à dire dans une structure locale similaire. Toutefois, comme le contexte d'une question est insuffisant pour effectuer une analyse sémantique congrue, cette structure locale est légère, c'est-à-dire qu'elle est limitée aux analyses morphologique et syntaxique, excluant donc la désambiguïsation sémantique et la phase d'expansion.

Certaines particularités doivent pourtant être signalées dans la conception de cette structure légère de la question. En effet, une grammaire particulière est mise en œuvre

« ...Domitien succéda à l'empereur Titus... »

Résultats avant expansion :

```
SUBJ(succéda[sn=1], Domitien[proper])
VARG[INDIR] (succéda[sn=1,], à, empereur[humain,])
NN(empereur, Titus[proper])
```

Structure de l'énoncé avec expansion :

```
SUBJ(succéda/remplacer, Domitien)
VARG[INDIR] (succéda, à, empereur)
VARG[DIR] (remplacer, empereur/chef/[...])
NN(empereur/chef/souverain/dots, Titus)

NMOD[INDIR] (successeur, de, empereur/chef/souverain/...)
NMOD(Domitien, successeur)
```

FIG. 7 – Construction de la structure informationnelle avec expansion

dans l'analyse de la question, qui permet deux adaptations de la structure. La première réside dans la relation **FOCUS**, qui permet de caractériser l'objet de la question, et donc la réponse attendue. Il s'agit d'une dépendance de marquage, qui identifie l'unité lexicale la plus significative de l'interrogation, c'est-à-dire la tête du groupe nominal lorsque l'interrogatif est un adjectif (« Qui est le beau-père de Galère ? » **FOCUS**(beau-père[PAR])) ou l'interrogatif lui-même si c'est un pronom (« Qui combattit les Parthes ? » **FOCUS**(qui[humain])). Elle permet d'identifier les traits sémantiques de l'objet de la question, et donc d'identifier la réponse dans les documents lorsque les autres éléments de la question se trouvent dans un énoncé de la base textuelle. Cette dépendance n'existant pas dans les documents – ni dans la structure informationnelle – étant donné ce qu'elle représente, elle devra ensuite être transmise comme un trait à l'intérieur de la structure légère, et supprimée comme dépendance cette structure légère.

La seconde adaptation de la structure locale à la question réside dans la suppression de toutes les informations ne relevant que du caractère interrogatif de cette question. Ainsi, l'interrogatif sera supprimé pour ne conserver, au sein des dépendances qui le contiennent, que les traits sémantiques qui lui sont propres et, le cas échéant, le trait **FOCUS**. Les dépendances purement fonctionnelles disparaissent également (dues au fonctionnement interne de XIP ou mettant en œuvre des mots-outils, des auxiliaires ou semi-auxiliaires), car elles ne sont pas porteuses d'information pertinente dans le cadre de cette application. La dépendance **FOCUS** sera éliminée de même, mais le trait subsiste dans les dépendances où doit apparaître le lexème sur lequel porte cette dépendance.

La recherche d'une réponse revient donc à mettre en correspondance la structure légère de la question, débarrassée de l'information propre à une interrogation, et des bribes de texte au travers de la structure informationnelle. Lorsqu'une information concordante à la structure légère est trouvée au sein de la même phrase dans la structure de l'information,

« De quel chef Domitien est-il le successeur ? »

Structure légère de la question :

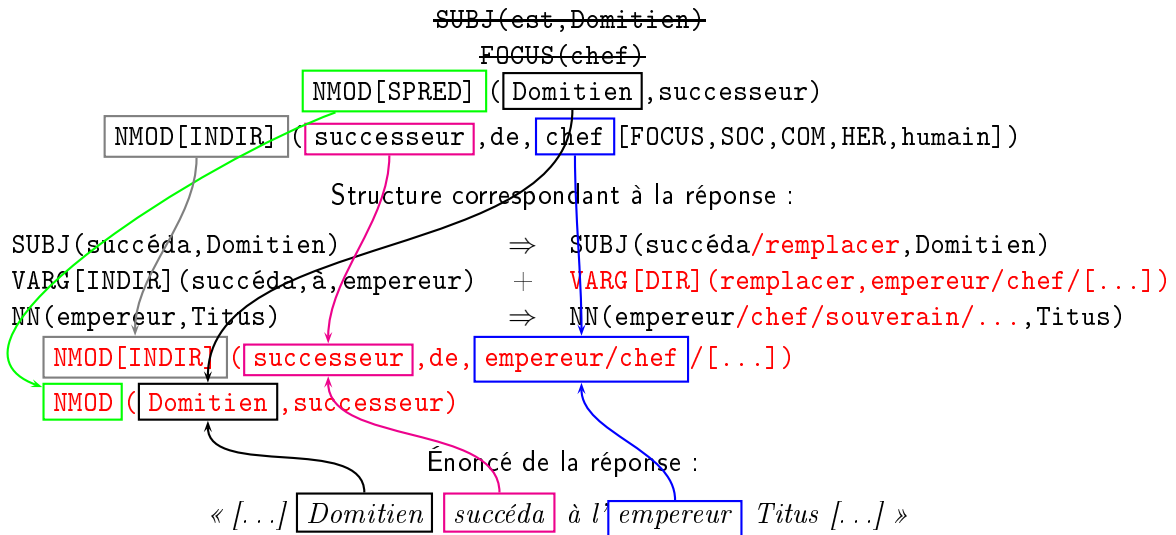


FIG. 8 – Exemple d'interrogation de la structure informationnelle avec son expansion

cette phrase est considérée comme une réponse pertinente à la question. Bien entendu, la réponse est considérée comme plus pertinente si une plus grande partie de l'information qui concorde est originale dans le texte, et moins pertinente à mesure que ces éléments concordants sont issus d'une expansion. La figure 8 page 11 illustre bien la mise en concordance d'une question avec sa réponse au travers de deux structures, l'une légère et purifiée, l'autre complète et enrichie d'expansions.

## 7 Conclusion

Nous avons présenté un système généraliste d'interrogation d'une base documentaire textuelle en langue naturelle. Ce système s'appuie sur des bases théoriques et sur des constatations pratiques pour proposer une méthode originale de structuration de l'information dans une base textuelle avec expansion des documents plutôt que des requêtes. L'ensemble des analyses et des enrichissements ont été effectués par des analyseurs linguistiques et les choix ont été réalisés suivant des indices contextuels et symbolistes issus de grammaires décrivant la langue.

Une analyse de ce système n'a pu être présentée ici faute de place. On peut en trouver le détail dans [Jacquemin, 2003]. Il montre la validité de la méthode – elle soutient la comparaison avec les meilleurs systèmes de sa catégorie dans la conférence TREC –, ainsi que certaines faiblesses, essentiellement liées à l'absence de résolution d'anaphores ou de hiérarchie sémantique. D'autre part, cette approche souffre, comme c'est habituel dans le domaine, de la représentation exclusivement lexicale de l'information, qui tient peu compte des mécanismes logiques. L'inférence, par exemple, n'est pas gérée actuellement, mais certaines approches statistique du lexique sont prometteuses à ce stade.

Par ailleurs, la présentation de ce système a été faite uniquement dans une optique de gestion de l'information. Cependant, il pourrait également se révéler un précieux outil d'étude de corpus écrit, dans la mesure où il peut être interrogé aisément et rapidement, que

tous les niveaux d'information linguistique sont disponibles à tout moment et qu'ils peuvent être individualisés sans problème. Ainsi, on peut facilement mêler dans une même requête des exigences lexicales, morphologiques, syntaxiques, sémantiques, de cooccurrence, et obtenir l'ensemble des réponses pertinentes quel que soit le corpus désiré, puisque ce système est automatique et qu'il accepte du texte tout venant avec une robustesse inhabituelle. Une telle approche à dominante linguistique semble dès lors se justifier, même si des améliorations peuvent et doivent y être apportées.

## Références

- [Aït-Mokhtar, 1998] AÏT-MOKHTAR, S. (1998). *L'analyse présyntaxique en une seule étape*. Thèse de doctorat, Université Clermont 2 Blaise Pascal, Clermont-Ferrand.
- [Bailly et Toro, 1947] BAILLY, R. et TORO, A. (1947). *Dictionnaire des synonymes de la langue française*. Larousse, Paris.
- [Brun et al., 2001] BRUN, C., JACQUEMIN, B. et SEGOND, F. (2001). Exploitation de dictionnaires électroniques pour la désambiguïsation sémantique lexicale. *Traitement Automatique des Langues*, 42(3):667–690.
- [Catherin, 1999] CATHERIN, L. (1999). The french wordnet. Rapport technique Deliverable 2D014, EuroWordNet.
- [Dubois et Dubois-Charlier, 1997] DUBOIS, J. et DUBOIS-CHARLIER, F. (1997). *Dictionnaire des verbes français*. Larousse, Paris. La première version de ce dictionnaire est électronique. Elle est accompagnée de son complément *Dictionnaire des mots*.
- [Ferret et al., 2002] FERRET, O., GRAU, B., HURAU-PLANTET, M., ILLOUZ, G., MONCEAUX, L., ROBBA, I. et VILNAT, A. (2002). Recherche de la réponse fondée sur la reconnaissance du focus de la question. In *Actes de TALN 2002*, pages 98–107.
- [Gaussier, 1999] GAUSSIER, r. (1999). Unsupervised learning of derivational morphology from inflectional lexicons. In *ACL'99 Workshop Proceedings on Unsupervised Learning in Natural Language Processing*, pages 24–30, College Park, Maryland, USA. ACL'99.
- [Harabagiu et al., 2000] HARABAGIU, S., MOLDOVAN, D., PSCA, M., MIHALCEA, R., SURDEANU, M., BUNESCU, R., GÎRJU, R., RUS, V. et MORARESCU, P. (2000). FALCON : Boosting knowledge for answer engines. In *Proceedings of Text REtrieval Conference*, pages 479–488.
- [Hull, 1999] HULL, D. A. (1999). Xerox trec-8 question answering track report. In *Proceedings of The Eighth Text Retrieval Conference*, pages 743–752. TREC-8.
- [Jacquemin, 2003] JACQUEMIN, B. (2003). *Construction et interrogation de la structure informationnelle d'une base documentaire en français*. Thèse de doctorat, Université Paris III Sorbonne Nouvelle, Paris. en cours.
- [Jacquemin, 2004a] JACQUEMIN, B. (2004a). Analyse et expansion des textes en question-réponse. In PURNELLE, G., FAIRON, C. et DISTER, A., éditeurs : *Le poids des mots. Actes des 7es journées internationales d'Analyse statistique des Données Textuelles*, volume 2, pages 633–641, Louvain-la-Neuve, Belgique. JADT04, Presses Universitaires de Louvain.
- [Jacquemin, 2004b] JACQUEMIN, B. (2004b). Dictionaries merger for text expansion in question answering. In ZOCK, M. et SAINT-DIZIER, P., éditeurs : *Proceedings of COLING. Enhancing and using electronic dictionaries*, Genève. COLING04.

- [Kando et Ishikawa, 2005] KANDO, N. et ISHIKAWA, H., éditeurs (2005). *Proceedings of the Fourth NTCIR Workshop on Research in Information Access Technologies Information Retrieval, Question Answering and Summarization*. NTCIR Workshop 4, National Institute of Informatics.
- [Peters *et al.*, 2005] PETERS, C., CLOUGH, P., GONZALO, J., JONES, G., KLUCK, M. et MAGNINI, B., éditeurs (2005). *Multilingual Information Access for Text, Speech and Images. 5th Workshop of the Cross-Language Evaluation Forum*, volume 3491, Bath, UK. CLEF 2004, Kluwer.
- [Reifler, 1955] REIFLER, E. (1955). The mechanical determination of meaning. In LOCKE, W. N. et BOOTH, A. D., éditeurs : *Machine translation of languages*, pages 136–164. John Wiley & Sons, New York.
- [Roux, 1999] ROUX, C. (1999). Phrase-driven parser. In *Proceedings of VEXTAL'99*, pages 235–240, Venezia, Italia. VEXTAL'99.
- [Roux et Jacquemin, 2002] ROUX, C. et JACQUEMIN, B. (2002). Storing and indexing of each level of the inner structure of a document with binary vector indexes, with the deepest level being the result of natural language processes. Déposition de proposition de brevet, Xerox.
- [Voorhees et Buckland, 2004] VOORHEES, E. M. et BUCKLAND, L. P., éditeurs (2004). *The Thirteenth Text REtrieval Conference Proceedings*, Gaithersburg, Maryland. TREC 2004, NIST-DARPA-ARDA.
- [Vossen, 1998] VOSSEN, P., éditeur (1998). *EuroWordNet : a multilingual database with lexical semantic networks*. Kluwer Academic Publishers, New York. Réédition de *Computer and the Humanities*, 32(2-3), 1998.
- [Weaver, 1949] WEAVER, W. (1949). Translation. In LOCKE, W. N. et BOOTH, A. D., éditeurs : *Machine translation of languages*, pages 15–23. John Wiley and Sons, New York.